

BANDWIDTH RESERVATION REUSE IN DYNAMICALLY ALLOCATED RING PROTECTION AND RESTORATION TECHNIQUE

Derek T. Mayweather

Jason C. Fan

Steven Gemelos

Robert F. Kalman

FIELD OF THE INVENTION

This invention relates to communication networks and, in particular, to networks employing rings.

BACKGROUND

As data services become increasingly mission-critical to businesses, service disruptions become increasingly costly. A type of service disruption that is of great concern is span outage, which may be due either to facility or equipment failures. Carriers of voice traffic have traditionally designed their networks to be robust in the case of facility outages, e.g. fiber breaks. As stated in the Telcordia GR-253 and GR-499 specifications for optical ring networks in the telecommunications infrastructure, voice or other protected services must not be disrupted for more than 60 milliseconds by a single facility outage. This includes up to 10 milliseconds for detection of a facility outage, and up to 50 milliseconds for rerouting of traffic.

A significant technology for implementing survivable networks meeting the above requirements has been SONET rings. A fundamental characteristic of such rings is that there are one (or more) independent physical links connecting adjacent nodes in the ring. Each link may be unidirectional, e.g. allow traffic to pass in a single direction, or may be bi-directional. A node is defined as a point where traffic

can enter or exit the ring. A single span connects two adjacent nodes, where a span consists of all links directly connecting the nodes. A span is typically implemented as either a two fiber or four fiber connection between the two nodes. In the two fiber case, each link is bi-directional, with half the traffic in each fiber going in the “clockwise” direction (or direction 0), and the other half going in the “counterclockwise” direction (or direction 1 opposite to direction 0). In the four fiber case, each link is unidirectional, with two fibers carrying traffic in direction 0 and two fibers carrying traffic in direction 1. This enables a communication path between any pair of nodes to be maintained on a single direction around the ring when the physical span between any single pair of nodes is lost. In the remainder of this document, references will be made only to direction 0 and direction 1 for generality.

There are 2 major types of SONET rings: unidirectional path-switched rings (UPSR) and bi-directional line-switched rings (BLSR). In the case of UPSR, robust ring operation is achieved by sending data in both directions around the ring for all inter-node traffic on the ring. This is shown in Fig. 1. This figure shows an N-node ring made up of nodes (networking devices) numbered from node 0 to node N-1 and interconnected by spans. In this document, nodes are numbered in ascending order in direction 0 starting from 0 for notational convenience. A link passing traffic from node i to node j is denoted by d_{ij} . A span is denoted by s_{ij} , which is equivalent to s_{ji} . In this document, the term span will be used for general discussion. The term link will be used only when necessary for precision. In this diagram, traffic from node 0 to node 5 is shown taking physical routes (bold arrows) in both direction 0 and direction 1. (In this document, nodes will be numbered sequentially in an increasing fashion in direction 0 for convenience. Node 0 will be used for examples.) At the receiving end, a special receiver implements “tail-end switching,” in which the receiver selects the data from one of the directions around the ring. The receiver can make this choice based on various performance monitoring (PM) mechanisms supported by SONET. This protection mechanism has the advantage that it is very simple, because no ring-level messaging is required to communicate a

span break to the nodes on the ring. Rather, the PM facilities built into SONET ensure that a “bad” span does not impact physical connectivity between nodes, since no data whatsoever is lost due to a single span failure.

Unfortunately, there is a high price to be paid for this protection. Depending on the traffic pattern on the ring, UPSR requires 100% extra capacity (for a single “hubbed” pattern) to 300% extra capacity (for a uniform “meshed” pattern) to as much as $(N-1)*100\%$ extra capacity (for an N node ring with a nearest neighbor pattern, such as that shown in Fig. 1) to be set aside for protection.

In the case of two-fiber BLSR, shown in Fig. 2A, data from any given node to another typically travels in one direction (solid arrows) around the ring. Data communication is shown between nodes 0 and 5. Half the capacity of each ring is reserved to protect against span failures on the other ring. The dashed arrows illustrate a ring that is typically not used for traffic between nodes 0 and 5 except in the case of a span failure or in the case of unusual traffic congestion.

In Fig. 2B, the span between nodes 6 and 7 has experienced a fault. Protection switching is now provided by reversing the direction of the signal from node 0 when it encounters the failed span and using excess ring capacity to route the signal to node 5. This switching, which takes place at the same nodes that detect the fault, is very rapid and is designed to meet the 50 millisecond requirement.

BLSR protection requires 100% extra capacity over that which would be required for an unprotected ring, since the equivalent of the bandwidth of one full ring is not used except in the event of a span failure. Unlike UPSR, BLSR requires ring-level signaling between nodes to communicate information on span cuts and proper coordination of nodes to initiate ring protection.

Though these SONET ring protection technologies have proven themselves to be robust, they are extremely wasteful of capacity. Additionally, both UPSR and BLSR depend intimately on the capabilities provided by SONET for their operation, and therefore cannot be readily mapped onto non-SONET transport mechanisms.

What is needed is a protection technology where no extra network capacity is consumed during “normal” operation (i.e., when all ring spans are operational), which is less tightly linked to a specific transport protocol, and which is designed to meet the Telcordia 50 millisecond switching requirement.

5 SUMMARY

A network protection and restoration technique and bandwidth reservation method is described that efficiently utilizes the total bandwidth in the network to overcome the drawbacks of the previously described networks, that is not linked to a specific transport protocol such as SONET, and that is designed to meet the Telcordia 50 millisecond switching requirement. The disclosed network includes two rings, wherein a first ring transmits data in a “clockwise” direction (or direction 0), and the other ring transmits data in a “counterclockwise” direction (or direction 1 opposite to direction 0). Additional rings may also be used. The traffic is removed from the ring by the destination node.

During normal operations (i.e., all spans operational and undegraded), data between nodes flows on the ring that provides the lowest-cost path to the destination node. If traffic usage is uniformly distributed throughout the network, the lowest-cost path is typically the minimum number of hops to the destination node. Thus, both rings are fully utilized during normal operations. Each node determines the lowest-cost path from it to every other node on the ring. To do this, each node must know the network topology.

A node monitors the status of each link for which it is at the receiving end, e.g. each of its ingress links, to detect a fault. The detection of such a fault causes a highest-priority link status broadcast message to be sent to all nodes. Processing at each node of the information contained in the link status broadcast message results in reconfiguration of a routing table within each node so as to identify the optimum routing of source traffic to the destination node after the fault. Hence, all nodes know the status of the network and all independently identify the optimal routing

path to each destination node when there is a fault in any of the links. The processing is designed to be extremely efficient to maximize switching speed.

Optionally, if it is desired to further increase the switching speed, an interim step can be used. A node that detects a link fault notifies its neighbor on the other
5 side of that span that a link has failed. Any node that detects an ingress link failure or that receives such a notification wraps inbound traffic headed for that span around onto the other ring. Traffic will be wrapped around only temporarily until the previously described rerouting of traffic is completed.

Since the remaining links will now see more data traffic due to the failed
10 link, traffic designated as "unprotected" traffic is given lower priority and may be dropped or delayed in favor of the "protected" traffic. Specific techniques are described for guaranteeing bandwidth availability for working and single failure traffic configurations, identifying a failed link, communicating the failed link to the other nodes, differentiating between protected and unprotected classes of traffic, and
15 updating the routing tables. Although the embodiments described transmit packets of data, the invention may be applied to any network transmitting frames, cells, or using any other protocol. Frames and cells are similar to packets in that all contain data and control information pertaining at least to the source and destination for the data. A single frame may contain multiple packets, depending on the protocol. A
20 cell may be fixed-size, depending on the protocol.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 illustrates inter-node physical routes taken by traffic from node 0 to node 5 using SONET UPSR, where a failure of spans between any single pair of nodes brings down only one of the two distinct physical routes for the traffic.

25 Fig. 2A illustrates an inter-node physical route taken by traffic from node 0 to node 5 using SONET two-fiber BLSR. Half of the capacity of each ring is reserved for protection, and half is used to carry regular traffic. The ring represented

with dashed lines is the ring in which protection capacity is used to reroute traffic due to the span failure shown.

Fig. 2B illustrates the bi-directional path taken by traffic from node 0 to node 5 using the SONET BLSR structure of Fig. 2A when there is a failure in the link
5 between nodes 6 and 7. Traffic is turned around when it encounters a failed link.

Fig. 3 illustrates a network in accordance with one embodiment of the present invention and, in particular, illustrates an inter-node physical route taken by traffic from node 0 to node 5.

Fig. 4 illustrates the network of Fig. 3 after a failure has occurred on the span
10 between nodes 6 and 7. When a failure occurs impacting a link or span on the initial path (e.g., between nodes 0 and 5), the traffic is rerouted at the ingress node to travel in the other direction around the ring to reach the destination node.

Fig. 5 illustrates the optional interim state of the network (based on wrapping traffic from one ring to the other) between that shown in Fig. 3 and that shown in
15 Fig. 4.

Fig. 6 illustrates pertinent hardware used in a single node.

Fig. 7 provides additional detail of the switching card and ring interface card in Fig. 6.

Fig. 8 is a flowchart illustrating steps used to identify a change in the status
20 of the network and to re-route traffic through the network.

Fig. 9 illustrates additional detail of the shelf controller card shown in Fig. 6.

DETAILED DESCRIPTION OF THE EMBODIMENTS

The purpose of the invention described herein is to achieve fast protection in a ring network while providing for efficient network capacity utilization. Certain
25 aspects of the preferred embodiment are:

- a. Transmission of a given packet between two nodes in only one direction around the ring (rather than in both directions as is done in SONET UPSR).
- b. Differentiation between “protected” and “unprotected” traffic classes.
- 5 c. A fast topology communication mechanism to rapidly communicate information about a span break to all nodes in the ring.
- d. A fast re-routing/routing table update mechanism to re-route paths impacted by a span break the other direction around the ring.
- e. An optional interim wrapping mechanism that may be used to further
10 increase protection switching speed.

These aspects are described in more detail below.

Unidirectional Transmission

A given packet/flow between two nodes is transmitted in only a single direction around the network (even when there is a span fault) and is removed from
15 the ring by the destination node, as is shown in Fig. 3 where node 0 transmits information to node 5 in only the direction indicated by the thick arrows. A transmission from node 5 to node 0 would only go through nodes 6 and 7 in the opposite direction. This allows for optimized ring capacity utilization since no capacity is set aside for protection.

20 The least-cost physical route is typically used for protected traffic. This is often the shortest-hop physical route. For example, a transmission from node 0 to node 2 would typically be transmitted via node 1. The shortest-hop physical route corresponds to the least-cost route when traffic conditions throughout the network are relatively uniform. If traffic conditions are not uniform, the least-cost physical
25 route from node 0 to node 2 can instead be the long path around the ring.

The removal of packets from the ring by the destination node ensures that traffic does not use more capacity than is necessary to deliver it to the destination node, thus enabling increased ring capacity through spatial reuse of capacity. An example of spatial reuse is the following. If 20% of span capacity is used up for traffic flowing from node 0 to node 2 via node 1, then the removal of this traffic from the ring at node 2 means that the 20% of span capacity is now available for any traffic flowing on any of the other spans in the ring (between nodes 2 and 3, nodes 3 and 4, etc.)

Protected and Unprotected Traffic Classes

In the case of unidirectional transmission described above, the loss of any span in the ring will result in a reduction in network capacity. This follows from the fact that traffic that would flow along a given span during normal operations must share the capacity of other spans in the case of a failure of that span. For example, Fig. 4 shows a span break between nodes 6 and 7. In contrast to Fig. 3, a transmission from node 0 to node 5 must now travel in a clockwise direction on another ring (illustrated by the thick arrows), adding to the traffic on that ring.

Because some network capacity is lost in the case of a span outage, a heavily loaded network with no capacity set aside for protection must suffer some kind of performance degradation as a result of such an outage. If traffic is classified into a “protected” class and an “unprotected” class, network provisioning and control can be implemented such that protected traffic service is unaffected by the span outage. This control is achieved through the use of bandwidth reservation management that processes provisioning requests considering the impact of a protection switch. In such a case, all of the performance degradation is “absorbed” by the unprotected traffic class via a reduction in average, peak, and burst bandwidth allocated to unprotected traffic on remaining available spans so that there is sufficient network capacity to carry all protected traffic. Traffic within the unprotected class can be further differentiated into various subclasses such that certain subclasses suffer more degradation than do others.

Fast Topology Communication Mechanism

Due to Telcordia requirements previously mentioned, the loss of a span in a ring must be rapidly sensed and communicated to all nodes in a ring.

In the case of a span outage, the node on the receiving end of each link within the span detects that each individual link has failed. If only a single link is out, then only the loss of that link is reported. Depending on the performance monitoring (PM) features supported by the particular communications protocol stack being employed, this detection may be based on loss of optical (or electrical) signal, bit error rate (BER) degradation, loss of frame, or other indications.

Each link outage must then be communicated to the other nodes. This is most efficiently done through a broadcast (store-and-forward) message (packet), though it could also be done through a unicast message from the detecting node to each of the other nodes in the network. This message must at least be sent out on the direction opposite to that leading to the broken span. The message must contain information indicating which link has failed.

Fast Source Node Re-routing Mechanism

When a link outage message is received by a given node, the node must take measures to re-route traffic that normally passed through the link. A possible sequence of actions is:

- a. Receive link outage message;
- b. Evaluate all possible inter-node physical routes (there are $2*(N-1)$ of them in an N node ring) to determine which ones are impacted by the loss of the link;
- c. Update routing tables to force all impacted traffic to be routed the other way around the ring; and

d. Update capacity allocated to unprotected traffic classes to account for reduced network capacity associated with the link outage. Details of how this capacity allocation is accomplished are not covered in this specification.

Being able to perform the operations above quickly requires that the various tables be properly organized to rapidly allow affected paths to be identified. Additionally, updates must be based either on computationally simple algorithms or on pre-calculated lookup tables.

Optional Interim Wrapping Mechanism

To increase the speed of protection switching, it may be desirable to take direct action at the node(s) detecting the fault, rather than waiting for re-routing to take place at all nodes. A possible sequence of actions is:

a. Upon detection of an ingress link fault, a node must transmit a neighbor fault notification message to the node on the other side of the faulty link. This notification is only required if there is a single link failure, as the node using the failed link as an egress link would not be able to detect that it had become faulty. In the event that a full span is broken, the failure to receive these notifications do not affect the following steps.

b. Upon detection of an ingress link fault or upon receipt of a neighbor fault notification message, a node must wrap traffic bound for the corresponding egress link on that span onto the other ring. This is shown in Fig. 5. Traffic from node 0 bound for node 5 is wrapped by node 7 onto the opposite ring because the span connecting node 7 to node 6 is broken.

The above steps are optional and should only be used if increased protection switching speed using this approach is required. This is because wrapping traffic from one ring onto the other uses up significantly more ring capacity than the standard approach described in this document. During the period, albeit short, between the start of wrapping and the completion of rerouting at source nodes, the

capacity that must be reserved for protection is as much as that required in two-fiber BLSR.

Specific Algorithms

Bandwidth Reservation for Protected and Unprotected Traffic Provisioning

5 This section describes the mechanism used to account for provisioned bandwidth on the ring. Define $C_{new}(j, k, 0)$ as a new simplex connection from node j to node k on ring 0 (the clockwise ring as shown in Fig. 3). Assume that $k > j$. If not, the representative node numbering across the ring (for this example) can be re-done so that $j=0$ and $k=k-j$. Similarly, $C_{new}(k, j, 1)$ would be a new simplex connection
10 from node k to node j on ring 1 (the counter-clockwise ring as shown in Fig. 3). Connection $C_{new}(j, k, 0)$ has a peak provisioned, or allowable, bandwidth of B . A connection may be provisioned either simplex or full-duplex, where a full-duplex connection consists of both $C_{new}(j, k, 0)$ and $C_{new}(k, j, 1)$ and accounting would be required for each direction. A given connection $C_{new}(j, k, 0)$ can be provisioned
15 as either transporting protected traffic or unprotected traffic.

Each link has a maximum traffic capacity of L . To determine if the link is full, all traffic on the link must be summed. The traffic may be broken into different categories. For example, if the bandwidth constraints for the ring are class-based (or other categories), the request must also contain the associated class (category).
20 Also, it is important to note that the provisioned traffic of each type may be weighted, but is nominally one. Further, for bursty traffic, peak bandwidth considerations should be made in the bandwidth accounting. For example, if three classes are supported (EF, AF, and BE), the amount of traffic per class that is allowed on a link can be governed through class-specific over-subscription
25 parameters c^{EF} , c^{AF} , c^{BE} as defined by

$$L \geq c^{EF} \cdot S^{EF} + c^{AF} \cdot S^{AF} + c^{BE} \cdot S^{BE}$$

where L is the high-speed link data rate and S is the aggregate traffic

Traffic matrices are used to determine the traffic provisioned in the ring. The elements of the matrices represent the aggregate bandwidth from a source node to a destination node. Thus the matrix element in row j and column k represents the aggregate bandwidth from node j to node k . There are two basic matrices defined:

5 **P** is the working traffic matrix for traffic requiring protection. The matrix element $\mathbf{P}[j, k]$ is the aggregate bandwidth from node j to node k of protected traffic. When a new wire is provisioned/removed, with protection, from node j to node k , with bandwidth B , B is added/subtracted to/from $\mathbf{P}[j, k]$. If a full-duplex wire is provisioned/removed, B is added/subtracted also to/from $\mathbf{P}[k, j]$.

10 **U** is the working traffic matrix for traffic not requiring protection. The matrix element $\mathbf{U}[j, k]$ is the aggregate bandwidth from node j to node k of unprotected traffic. When a new wire is provisioned/removed, without protection, from node j to node k , with bandwidth B , B is added/subtracted to/from $\mathbf{U}[j, k]$. If a full-duplex wire is provisioned/removed, B is added/subtracted also to/from $\mathbf{U}[k, j]$.

15 The traffic flow around the ring is bi-directional. Both clockwise and counterclockwise rings carry traffic. Clockwise and counter-clockwise rings will have its own set of basic traffic matrices. For a class-based category system, for EF traffic in the clockwise direction, there are \mathbf{P}_C^{EF} and \mathbf{U}_C^{EF} and for the counter-clockwise direction there are \mathbf{P}_{CC}^{EF} and \mathbf{U}_{CC}^{EF} .

20 Using the construct above, several checks can be made to determine if the bandwidth is available to support a new connection. These checks include verifying the bandwidth is available to support the working traffic configuration and every possible fault traffic configuration.

25 Using the constructs above, if $\text{Cnew}(j, k, 0)$ is provisioned, B is added to the $\mathbf{P}_C[j, k]$ element in the population matrix. Then the following class-based category span loading algorithm is run to verify the bandwidth on each span is available for the working configuration.

```

for (x=0 to N-1) { //spans 0 to N-1 for an N node network//
    ScEF[x] = 0; //Span X utilization due to EF traffic
    ScAF[x] = 0; //Span X utilization due to AF traffic
    ScBE[x] = 0; //Span X utilization due to BE traffic

```

5

```

    for (j = (1+x) to (N+x) ) {
        for (k = (1+x) to j ) {
            ScEF[x] = ScEF[x] + PCEF(j mod N, k mod N);
            ScEF[x] = ScEF[x] + UCEF(j mod N, k mod N);

```

10

```

            ScAF[x] = ScAF[x] + PCAF(j mod N, k mod N);
            ScAF[x] = ScAF[x] + UCAF(j mod N, k mod N);

```

15

```

            ScBE[x] = ScBE[x] + PCBE(j mod N, k mod N);
            ScBE[x] = ScBE[x] + UCBE(j mod N, k mod N);
        }
    }

```

```

    Sc[x] = cEF*ScEF[x] + cAF*ScAF[x] + cBE*ScBE[x];

```

20

```

    //Total Span X Utilization//

```

```

    if (Sc[x] > L) reject_provisioning_attempt=1;
}

```

25

If a rejection indication is not provided to the higher layer, the single failure configurations must be checked. To develop a single failure configuration, one-by one, a single link, w , is failed, where w is between node w and node $w+1$ on the clock wise ring. The traffic matrices are populated as discussed above; however, traffic that traversed link w is now switched at the source to the other ring. For each

30

```

    if (k >= j) {
        if (w>=k or w<j))

```

```

        Add crossconnect bandwidth to Pc[j,k];
    Else
        Add crossconnect bandwidth to Pcc[j,k];
    }
5   else {
        if (w>=j or w<k))
            Add crossconnect bandwidth to Pcc[j,k];
        Else
            Add crossconnect bandwidth to Pc[j,k];
10

```

For crossconnect C(j,k,1), the matrix is populated as follows:

```

15   if (k >= j) {
        if (w>=j or w<k))
            Add crossconnect bandwidth to Pcc[j,k];
        Else
            Add crossconnect bandwidth to Pc[j,k];
    }
20   else {
        if (w<=j or w>k))
            Add crossconnect bandwidth to Pc[j,k];
        Else
            Add crossconnect bandwidth to Pcc[j,k];
25

```

The unprotected crossconnects are provisioned as before, independent of the single failed link.

Once the single failure traffic configuration is generated as described, the same span loading algorithm described above is computed. Based upon the result,

30 the reject or accept indication is provided to the higher layer. This is performed for

each link in the clockwise and counter-clockwise direction. A failure of node N corresponds to a failure of links between nodes N-1 and N+1.

Fast Topology Communication Mechanism

5 This section describes a specific fast mechanism for communicating topology changes to the nodes in a ring network. The mechanism for communicating information about a span or link break or degradation from a node to all other nodes on a ring is as follows.

10 A link status message is sent from each node detecting any link break or degradation on ingress links to the node, e.g. links for which the node is on the receiving end. (Therefore, for a single span break the two nodes on the ends of the span will each send out a link status message reporting on the failure of a single distinct ingress link.) This message may be sent on the ring direction opposite the link break or on both ring directions. For robustness, it is desirable to send the message on both ring directions. In a network that does not wrap messages from one ring direction to the other ring direction, it is required that the message be sent on both ring directions to handle failure scenarios such as that in Fig. 4. The message may also be a broadcast or a unicast message to each node on the ring. For robustness and for capacity savings, it is desirable to use broadcast. In particular, 15 broadcast ensures that knowledge of the link break will reach all nodes, even those that are new to the ring and whose presence may not be known to the node sending the message. In either case, the mechanism ensures that the propagation time required for the message to reach all nodes on the ring is upper bounded by the time required for a highest priority message to travel the entire circumference of the ring. 20 It is desirable that each mechanism also ensure that messages passing through each node are processed in the fastest possible manner. This minimizes the time for the message to reach all nodes in the ring. 25

The link status message sent out by a node should contain at least the following information: source node address, link identification of the broken or degraded link for which the node is on the receive end, and link status for that link. For simplicity of implementation, the link status message can be expanded to

5 contain link identification and status for all links for which the node is on the receive end. The link identification for each link, in general, should contain at least the node address of the node on the other end of the link from the source node and the corresponding physical interface identifier of the link's connection to the destination node. The mechanism by which the source node obtains this information is found in

10 the co-pending application entitled "Dual-Mode Virtual Network Addressing," Serial No. _____, filed herewith by Jason Fan et al., assigned to the present assignee and incorporated herein by reference. The physical interface identifier is important, for example, in a two-node network where the address of the other node is not enough to resolve which link is actually broken or degraded. Link status

15 should indicate the level of degradation of the link, typically expressed in terms of measured bit error rate on the link (or in the event that the link is broken, a special identifier such as 1).

The link status message may optionally contain two values of link status for each link in the event that protection switching is non-revertive. An example of

20 non-revertive switching is illustrated by a link degrading due to, for example, temporary loss of optical power, then coming back up. The loss of optical power would cause other nodes in the network to protection switch. The return of optical power, however, would not cause the nodes to switch back to default routes in the case of non-revertive switching until explicitly commanded by an external

25 management system. The two values of link status for each link, therefore, may consist of a status that reflects the latest measured status of the link (previously described) and a status that reflects the worst measured status (or highest link cost) of the link since the last time the value was cleared by an external management system.

The link status message can optionally be acknowledged by the other nodes. In the event that the message is not acknowledged, it must be sent out multiple times to ensure that it is received by all other nodes. In the event that the message requires acknowledgement on receipt, it must be acknowledged by all expected recipient
5 nodes within some time threshold. If not, the source node may choose to re-send the link status message to all expected recipients, or re-send the link status message specifically to expected recipients that did not acknowledge receipt of the message.

Fast Source Node Re-routing Mechanism

This section describes a mechanism which allows a node in a ring network to
10 rapidly re-route paths that cross broken links. The following describes a fast source node re-routing mechanism when node 0 is the source node.

For each destination node j , a cost is assigned to each output direction (0 and 1) from node 0 on the ring. A preferred direction for traffic from nodes 0 to j is selected based on the direction with the lowest cost. For simplicity, the mechanism
15 for reassigning costs to the path to each destination node for each output direction from node 0 operates with a constant number of operations, irrespective of the current condition of the ring. (The mechanism may be further optimized to always use the minimum possible number of operations, but this will add complexity to the algorithm without significantly increasing overall protection switching speed.) The
20 mechanism for reassigning an output direction to traffic packets destined for a given node based on the path cost minimizes the time required to complete this reassignment.

A table is maintained at each node with the columns Destination Node, direction 0 cost, and direction 1 cost. An example is shown as Table 1. The
25 computation of the cost on a direction from node 0 (assuming node 0 as the source) to node j may take into account a variety of factors, including the number of hops from source to destination in that direction, the cumulative normalized bit error rate from source to destination in that direction, and the level of traffic congestion in that

direction. Based on these costs, the preferred output direction for traffic from the source to any destination can be selected directly. The example given below assumes that the costs correspond only to the normalized bit error rate from source to destination in each direction. The cost on a given link is set to 1 if the measured bit error rate is lower than the operational bit error rate threshold. Conveniently, if all links are fully operational, the cumulative cost from node 0 to node j will be equal to the number of hops from node 0 to node j if there is no traffic congestion. Traffic congestion is not taken into account in this example.

For a representative ring with a total of 8 nodes (in clockwise order 0, 1, 2, 3, 4, 5, 6, 7), the table's normal operational setting at node 0 is:

Table 1. Preferred direction table at node 0

Destination Node	Direction 0 cost	Direction 1 cost	Preferred Direction
1	1	7	0
2	2	6	0
3	3	5	0
4	4	4	0
5	5	3	1
6	6	2	1
7	7	1	1

The preferred direction is that with the lower cost to reach destination node j. In the event that the costs to reach node j on direction 0 and on direction 1 are equal, then either direction can be selected. (Direction 0 is selected in this example.) The normal operational cost for each physical route (source to destination) is computed from the link status table shown in Table 3.

The pseudocode for selection of the preferred direction is:

For j=1 to N-1 {N is the total number of nodes in the ring}

Update direction 0 cost ($\text{dir_0_cost}(j)$) and direction 1 cost ($\text{dir_1_cost}(j)$) for each destination node j ; {expanded later in this section}

{HYST_FACT is the hysteresis factor to prevent a ping-pong effect due to BER variations in revertive networks. A default value for this used in SONET is 10}

5 If ($\text{dir_0_cost}(j) < \text{dir_1_cost}(j)/\text{HYST_FACT}$),
 $\text{dir_preferred}(j) = 0$;
 Else if ($\text{dir_1_cost}(j) < \text{dir_0_cost}(j)/\text{HYST_FACT}$),
 $\text{dir_preferred}(j) = 1$;
 Else if $\text{dir_preferred}(j)$ has a pre-defined value,
 {This indicates that $\text{dir_preferred}(j)$ has been previously set to a preferred direction and thus should not change if the above two conditions were not met}
 $\text{dir_preferred}(j)$ does not change;
 Else if $\text{dir_preferred}(j)$ does not have a pre-defined value,
 if $\text{dir_0_cost}(j) < \text{dir_1_cost}(j)$,
 $\text{dir_preferred}(j) = 0$;
 Else if $\text{dir_1_cost}(j) < \text{dir_0_cost}(j)$,
 $\text{dir_preferred}(j) = 1$;
 Else
 $\text{dir_preferred}(j) = 0$;
 End {else if $\text{dir_preferred}(j)$ does not have a pre-defined value}
 End {for loop j }

25 The link status table (accessed by a CPU at each node) is used to compute the costs in the preferred direction table above. The link status table's normal operational setting looks like:

Table 3. Link status table (identical at every node)

Link Identifier, direction 0	Link Identifier, direction 1	Direction 0 cost	Direction 1 cost
d ₀₁	d ₁₀	1	1
d ₁₂	d ₂₁	1	1
d ₂₃	d ₃₂	1	1
d ₃₄	d ₄₃	1	1
d ₄₅	d ₅₄	1	1
d ₅₆	d ₆₅	1	1
d ₆₇	d ₇₆	1	1
d ₇₀	d ₀₇	1	1

The cost for each link d_{ij} is the normalized bit error rate, where the measured bit error rate on each link is divided by the default operational bit error rate (normally 10E-9 or lower). In the event that the normalized bit error rate is less than 1 for a link, the value entered in the table for that link is 1.

The pseudocode for the line “Update direction 0 cost and direction 1 cost” for each node j in the pseudocode for selection of preferred direction uses the link status table shown in Table 3 as follows:

{Initialization of Linkcostsum values in each direction. These variables are operated on inside the for loop below to generate dir_0_cost(j) and dir_1_cost(j).}

Linkcostsum_{dir 0} = 0;

{ Linkcostsum_{dir 1} is the sum of link costs all the way around the ring in direction 1, starting at node 0 and ending at node 0. }

Linkcostsum_{dir 1} = sum over all links(Linkcost_{dir 1});

For $j=0$ to $N-1$ { N is the total number of nodes in the ring}

{MAX_COST is the largest allowable cost in the preferred direction table. Linkcost_{dir 0, link i,j} is the cost of the link in direction 0 from node i to node j .}

```

If (Linkcostsumdir 0 < MAX_COST)
    Linkcostsumdir 0 = Linkcostsumdir 0 + Linkcostdir 0, link j, (j+1)
    modN;
else
5     Linkcostsumdir 0 = MAX_COST;
    dir_0_cost(j) = Linkcostsumdir 0;
    If (Linkcostsumdir 1 < MAX_COST)
        Linkcostsumdir 1 = Linkcostsumdir 1 - Linkcostdir 1, link (j+1) modN,
        j;
10    else
        Linkcostsumdir 1 = MAX_COST;
        dir_1_cost(j) = Linkcostsumdir 1;
    End {for loop j}

```

15 The update of the link status table is based on the following pseudocode:

{This version of the pseudocode assumes more than 2 nodes in the ring}

If (linkstatusmessage.source = node i) and (linkstatusmessage.neighbor = node j) and (direction = 0)

Linkcost_{dir 0, link i, j} = linkstatusmessage.status;

20 else if (linkstatusmessage.source = node i) and
(linkstatusmessage.neighbor = node j) and (direction = 1) Linkcost_{dir 1, link j, i} =
linkstatusmessage.status;

25 In the event that a link is broken, the linkstatusmessage.status for that link is
a very large value. In the event that a link is degraded, the linkstatusmessage.status
for that link is the measured bit error rate on that link divided by the undegraded bit
error rate of that link. All undegraded links are assumed to have the same
undegraded bit error rate.

The link status table may optionally contain two cost columns per direction to handle non-revertive switching scenarios. These would be measured cost (equivalent to the columns currently shown in Table 3) and non-revertive cost. The non-revertive cost column for each direction contains the highest value of link cost reported since the last time the value was cleared by an external management system. This cost column (instead of the measured cost) would be used for preferred direction computation in the non-revertive switching scenario. The preferred direction table may also optionally contain two cost columns per direction, just like the link status table. It may also contain two preferred direction columns, one based on the measured costs and the other based on the non-revertive costs. Again, the non-revertive cost columns would be used for computations in the non-revertive switching scenario.

As an example, assume that the clockwise link between node 2 and node 3 is degraded with factor a (where $a > \text{HYST_FACT}$), the clockwise link between node 4 and node 5 is broken (factor MAX), the counterclockwise link between node 1 and node 2 is degraded with factor b (where $b > \text{HYST_FACT}$), and the counterclockwise link between node 5 and node 6 is degraded with factor c (where $c < a/\text{HYST_FACT}$). The link status table for this example is shown in Table 5.

Table 5. Example of link status table with degraded and broken links

Link Identifier, direction 0	Link Identifier, direction 1	Direction 0 cost (clockwise)	Direction 1 cost (counterclockwise)
d ₀₁	d ₁₀	1	1
d ₁₂	d ₂₁	1	b
d ₂₃	d ₃₂	a	1
d ₃₄	d ₄₃	1	1
d ₄₅	d ₅₄	MAX	1
d ₅₆	d ₆₅	1	c
d ₆₇	d ₇₆	1	1
d ₇₀	d ₀₇	1	1

The costs of the links needed between the source node and destination node are added to determine the total cost.

5 The preferred direction table for the source node 0 is then:

Table 7. Example of preferred direction table with degraded and broken links

Destination Node	Direction 0 cost (clockwise)	Direction 1 cost (counterclockwise)	Preferred Direction
1	1	c+b+5	0
2	2	c+5	0
3	a+2	c+4	1
4	a+3	c+3	1
5	MAX	c+2	1
6	MAX	2	1
7	MAX	1	1

(In the selection of the preferred direction, it is assumed that HYST_FACT = 10.)

Once these preferred directions are determined, a corresponding mapping table of destination node to preferred direction in packet processors on the data path is modified to match the above table.

Neighbor Fault Notification in Optional Interim Wrapping Mechanism

This section describes a specific fast mechanism for communication of a fault notification from the node on one side of the faulty span to the node on the other side. This mechanism, as described previously, is only necessary in the event of a single link failure, since the node using that link as its egress link cannot detect that it is faulty.

A neighbor fault notification message is sent from each node detecting any link break or degradation on an ingress link to the node. The message is sent on each egress link that is part of the same span as the faulty ingress link. To ensure that it is received, the notification message can be acknowledged via a transmission on both directions around the ring. If it is not acknowledged, then the transmitting node must send the notification multiple times to ensure that it is received. The message is highest priority to ensure that the time required to receive the message at the destination is minimized.

The neighbor fault notification message sent out by a node should contain at least the following information: source node address, link identification of the broken or degraded link for which the node is on the receive end, and link status for that link. For simplicity of implementation, the neighbor fault notification message may be equivalent to the link status message broadcast to all nodes that has been previously described.

Mechanisms to Provide Provisioning and Routing Information to Tributary Interface Cards

Fig. 9 illustrates one shelf controller card 62 in more detail. The shelf controller 62 obtains status information from the node and interfaces with a network management system. The shelf controller 62 both provisions other cards within the device 20 and obtains status information from the other cards. In addition, the shelf controller interfaces with an external network management system and with other types of external management interfaces. The software applications controlling these functions run on the CPU 92. The CPU may be an IBM/Motorola MPC750 microprocessor.

A memory 93 represents memories in the node. It should be understood that there may be distributed SSRAM, SDRAM, flash memory and EEPROM to provide the necessary speed and functional requirements of the system.

The CPU is connected to a PCI bridge 94 between the CPU and various types of external interfaces. The bridge may be an IBM CPC700 or any other suitable type.

Ethernet controllers 96 and 102 are connected to the PCI bus. The controller may be an Intel 21143 or any other suitable type.

An Ethernet switch 98 controls the Layer 2 communication between the shelf controller and other cards within the device. This communication is via control lines on the backplane. The layer 2 protocol used for the internal communication is 100BaseT switched Ethernet. This switch may be a Broadcom BCM5308 Ethernet switch or any other suitable type.

The output of the Ethernet switch must pass through the Ethernet Phy block 100 before going on the backplane. The Ethernet Phy may be a Bel Fuse, Inc., S558 or any other suitable type that interfaces directly with the Ethernet switch used.

The output of the Ethernet controller 102 must pass through an Ethernet Phy 104 before going out the network management system (NMS) 10/100 BaseT Ethernet port. The Ethernet Phy may be an AMD AM79874 or any other suitable type.

5 Information is delivered between applications running on the shelf controller CPU and applications running on the other cards via well-known mechanisms including remote procedure calls (RPCs) and event-based notification. Reliability is provided via TCP/IP or via UDP/IP with retransmissions.

10 Provisioning of cards and ports via an external management system is via the NMS Ethernet port. Using a well-known network management protocol such as the Simple Network Management Protocol (SNMP), the NMS can control a device via the placement of an SNMP agent application on the shelf controller CPU. The SNMP agent interfaces with a shelf manager application. The shelf manager application is primarily responsible for the provisioning on tributary interface cards in 52.

15 Communication from the shelf controller onto the ring is via the switching card CPU. This type of communication is important for sending SNMP messages to remote devices on the ring from an external management system physically connected to the shelf. The bandwidth management that determines whether provisioning is accepted runs on the shelf controller or an external workstation.

DESCRIPTION OF HARDWARE

25 Fig. 6 illustrates the pertinent functional blocks in each node. Node 0 is shown as an example. Each node is connected to adjacent nodes by ring interface cards 30 and 32. These ring interface cards convert the incoming optical signals on fiber optic cables 34 and 36 to electrical digital signals for application to switching card 38.

Fig. 7 illustrates one ring interface card 32 in more detail showing the optical transceiver 40. An additional switch in card 32 may be used to switch between two switching cards for added reliability. The optical transceiver may be a Gigabit Ethernet optical transceiver using a 1300 nm laser, commercially available.

5 The serial output of optical transceiver 40 is converted into a parallel group of bits by a serializer/deserializer (SERDES) 42 (Fig. 6). The SERDES 42, in one example, converts a series of 10 bits from the optical transceiver 40 to a parallel group of 8 bits using a table. The 10 bit codes selected to correspond to 8 bit codes meet balancing criteria on the number of 1's and 0's per code and the maximum
10 number of consecutive 1's and 0's for improved performance. For example, a large number of sequential logical 1's creates baseline wander, a shift in the long-term average voltage level used by the receiver as a threshold to differentiate between 1's and 0's. By utilizing a 10-bit word with a balanced number of 1's and 0's on the backplane, the baseline wander is greatly reduced, thus enabling better AC coupling
15 of the cards to the backplane.

When the SERDES 42 is receiving serial 10-bit data from the ring interface card 32, the SERDES 42 is able to detect whether there is an error in the 10-bit word if the word does not match one of the words in the table. The SERDES 42 then generates an error signal. The SERDES 42 uses the table to convert the 8-bit code
20 from the switching card 38 into a serial stream of 10 bits for further processing by the ring interface card 32. The SERDES 42 may be a model VSC 7216 by Vitesse or any other suitable type.

A media access controller (MAC) 44 counts the number of errors detected by the SERDES 42, and these errors are transmitted to the CPU 46 during an interrupt
25 or pursuant to polling mechanism. The CPU 46 may be a Motorola MPC860DT microprocessor. Later, it will be described what happens when the CPU 46 determines that the link has degraded sufficiently to take action to cause the nodes to re-route traffic to avoid the faulty link. The MAC 44 also removes any control words forwarded by the SERDES and provides OSI layer 2 (data-link) formatting

for a particular protocol by structuring a MAC frame. MACs are well known and are described in the book "Telecommunication System Engineering" by Roger Freeman, third edition, John Wiley & Sons, Inc., 1996, incorporated herein by reference in its entirety. The MAC 44 may be a field programmable gate array.

5 The packet processor 48 associates each of the bits transmitted by the MAC 44 with a packet field, such as the header field or the data field. The packet processor 48 then detects the header field of the packet structured by the MAC 44 and may modify information in the header for packets not destined for the node. Examples of suitable packet processors 48 include the XPIF-300 Gigabit Bitstream Processor or the EPIF 4-L3C1 Ethernet Port L3 Processor by MMC Networks, whose data sheets are incorporated herein by reference.

10 The packet processor 48 interfaces with an external search machine/memory 47 (a look-up table) that contains routing information to route the data to its intended destination. The updating of the routing table in memory 47 will be discussed in detail later.

15 A memory 49 in Fig. 6 represents all other memories in the node, although it should be understood that there may be distributed SSRAM, SDRAM, flash memory, and EEPROM to provide the necessary speed and functional requirements of the system

20 The packet processor 48 provides the packet to a port of the switch fabric 50, which then routes the packet to the appropriate port of the switch fabric 50 based on the packet header. If the destination address in the packet header corresponds to the address of node 0 (the node shown in Fig. 6), the switch fabric 50 then routes the packet to the appropriate port of the switch fabric 50 for receipt by the designated node 0 tributary interface card 52 (Fig. 5) (to be discussed in detail later). If the packet header indicates an address other than to node 0, the switch fabric 50 routes the packet through the appropriate ring interface card 30 or 32 (Fig. 5). Control packets are routed to CPU 46. Such switching fabrics and the routing techniques

used to determine the path that packets need to take through switch fabrics are well known and need not be described in detail.

One suitable packet switch is the MMC Networks model nP5400 Packet Switch Module, whose data sheet is incorporated herein by reference. In one
5 embodiment, four such switches are connected in each switching card for faster throughput. The switches provide packet buffering, multicast and broadcast capability, four classes of service priority, and scheduling based on strict priority or weighted fair queuing.

A packet processor 54 associated with one or more tributary interface cards,
10 for example, tributary interface card 52, receives a packet from switch fabric 50 destined for equipment (e.g., a LAN) associated with tributary interface card 52. Packet processor 54 is bi-directional, as is packet processor 48. Packet processors 54 and 48 may be the same model processors. Generally, packet processor 54 detects the direction of the data through packet processor 54 as well as accesses a
15 routing table memory 55 for determining some of the desired header fields and the optimal routing path for packets heading onto the ring, and the desired path through the switch for packets heading onto or off of the ring. This is discussed in more detail later. When the packet processor 54 receives a packet from switch fabric 50, it forwards the packet to a media access control (MAC) unit 56, which performs a
20 function similar to that of MAC 44, which then forwards the packet to the SERDES 58 for serializing the data. SERDES 58 is similar to SERDES 42.

The output of the SERDES 58 is then applied to a particular tributary interface card, such as tributary interface card 52 in Fig. 5, connected to a backplane 59. The tributary interface card may queue the data and route the data to a particular
25 output port of the tributary interface card 52. Such routing and queuing by the tributary interface cards may be conventional and need not be described in detail. The outputs of the tributary interface cards may be connected electrically, such as via copper cable, to any type of equipment, such as a telephone switch, a router, a LAN, or other equipment. The tributary interface cards may also convert electrical

signals to optical signals by the use of optical transceivers, in the event that the external interface is optical.

In one embodiment, the above-described hardware processes bits at a rate greater than 1Gbps.

5 Functions of Hardware During Span Failure/Degradation

Fig. 8 is a flow chart summarizing the actions performed by the network hardware during a span failure or degradation. Since conventional routing techniques and hardware are well known, this discussion will focus on the novel characteristics of the preferred embodiment.

10 In step 1 of Fig. 8, each of the nodes constantly or periodically tests its links with neighboring nodes. The MAC 44 in Fig. 7 counts errors in the data stream (as previously described) and communicates these errors to the CPU 46. The CPU compares the bit error rate to a predetermined threshold to determine whether the link is satisfactory. An optical link failure may also be communicated to the CPU.
15 CPU 46 may monitor ingress links from adjacent devices based on error counting by MAC 44 or based on the detection of a loss of optical power on ingress fiber 36. This detection is performed by a variety of commercially available optical transceivers such as the Lucent NetLight transceiver family. The loss of optical power condition can be reported to CPU 46 via direct signaling over the backplane
20 (such as via I2C lines), leading to an interrupt or low-level event at the CPU.

In step 2, the CPU 46 determines if there is a change in status of an adjacent link. This change in status may be a fault (bit error rate exceeding threshold) or that a previously faulty link has been repaired. It will be assumed for this example that node 6 sensed a fault in ingress link connecting it to node 7.

25 If there is no detection of a fault in step 2, no change is made to the network. It is assumed in Fig. 8 that adjacent nodes 6 and 7 both detect faults on ingress links connecting node 6 to node 7. The detection of a fault leads to an interrupt or low-

level event (generated by MAC 44) sent through switch fabric 50 to CPU 46 signaling the change in status.

In optional step 3, nodes 6 and 7 attempt to notify each other directly of the ingress link fault detected by each. The notification sent by node 6, for example, is sent on the egress link of node 6 connected to node 7. If the entire span is broken, these notifications clearly do not reach the destination. They are useful only if a single link within a span is broken. This is because a node has no way to detect a fiber break impacting an egress link. Based on this notification, each node can then directly wrap traffic in the fashion shown in Fig. 5. The wrapping of traffic in node 6 is performed through a configuration command from CPU 46 to packet processor 48 connected as shown in Fig. 7 to ring interface card 32 (assuming that links from ring interface card 32 connect to node 7). After receiving this command, packet processor 48 loops back traffic through the switching fabric and back out ring interface card 30 that it normally would send directly to node 7.

Each communication by a node of link status is associated with a session number. A new session number is generated by a node only when it senses a change in the status of a neighboring node. As long as the nodes receive packets with the current session number, then the nodes know that there is no change in the network. Both nodes 6 and 7 increment the session number stored at each node upon detection of a fault at each node.

In step 4, both node 6 and node 7 then broadcast a link status message, including the new session number, conveying the location of the fault to all the nodes. Each node, detecting the new session number, forwards the broadcast to its adjacent node.

A further description of the use of the session number in general topology reconfiguration scenarios, of which a link or span failure is one, is found in the co-pending application entitled "Dual-Mode Virtual Network Addressing," by Jason Fan et al., assigned to the present assignee and incorporated herein by reference.

In step 5, the identity of the fault is then used by the packet processor 54 in each node to update the routing table in memory 55. Routing tables in general are well known and associate a destination address in a header with a particular physical node to which to route the data associated with the header. Each routing table is then configured to minimize the cost from a source node to a destination node. Typically, if the previously optimized path to a destination node would have had to go through the faulty link, that route is then updated to be transmitted through the reverse direction through the ring to avoid the faulty route. The routing table for each of the packet processors 54 in each node would be changed as necessary depending upon the position of the node relative to the faulty link. Details of the routing tables have been previously described.

In one embodiment, each of the nodes must acknowledge the broadcast with the new session number, and the originating node keeps track of the acknowledgments. After a time limit has been exceeded without receiving all of the acknowledgments, the location of the fault is re-broadcast without incrementing the sequence number.

Accordingly, all nodes store the current topology of the ring, and all nodes may independently create the optimum routing table entries for the current configuration of the ring.

In step 6, the routing table for each node has been updated and data traffic resumes. Accordingly, data originating from a LAN connected to a tributary interface card 52 (Fig. 5) has appended to it an updated routing header by packet processor 54 for routing the data through switch fabric 50 to the appropriate output port for enabling the data to arrive at its intended destination. The destination may be the same node that originated the data and, thus, the switch fabric 50 would wrap the data back through a tributary interface card in the same node. Any routing techniques may be used since the invention is generally applicable to any protocol and routing techniques.

Since some traffic around the ring must be re-routed in order to avoid the faulty link, and the bandwidths of the links are fixed, the traffic to be transmitted around the healthy links may exceed the bandwidth of the healthy links.

Accordingly, some lower priority traffic may need to be dropped or delayed, as

5 identified in step 7. Generally, the traffic classified as “unprotected” is dropped or delayed as necessary to support the “protected” traffic due to the reduced bandwidth.

In one embodiment, the packet processor 54 detects the header that identifies the data as unprotected and drops the packet, as required, prior to the packet being applied to the switch fabric 50. Voice traffic is generally protected.

10 In step 8, switch fabric 50 routes any packet forwarded by packet processor 54 to the appropriate output port for transmission either back into the node or to an adjacent node.

15 The above description of the hardware used to implement one embodiment of the invention is sufficient for one of ordinary skill in the art to fabricate the invention since the general hardware for packet switching and routing is very well known. One skilled in the art could easily program the MACs, packet processors, CPU 46, and other functional units to carry out the steps describe herein. Firmware or software may be used to implement the steps described herein.

20 While particular embodiments of the present invention have been shown and described, it will be obvious to those skilled in the art that changes and modifications may be made without departing from this invention in its broader aspects and, therefore, the appended claims are to encompass within their scope all such changes and modifications as fall within the true spirit and scope of this invention.

25